# Social Environments Modeling From Kinect Data in Robotics Applications

Catarina Lima and João Silva Sequeira

*Instituto Superior Técnico / Institute for Systems and Robotics, University of Lisbon, Lisbon, Portugal*

Keywords:     Social Robots, Human-Robot Interaction, Environment Modeling, Kinect, Robot Behavior.

Abstract:     This paper addresses the modeling of social environments from range information obtained from a Kinect sensor. The modeling is restricted to events representing the existence of movement in front of the sensor. A deterministic model based on a power law and probabilistic models based in Weibull and Lognormal distributions are considered. Real experiments in a hospital ward are presented together with a discussion on the relevance of these models to improve the acceptance of social robots in non lab social environments.

## 1 INTRODUCTION

The importance of social interactions involving humans and robots is on the rise. Researchers want to build not only functional robots but also social-functional robots which can help and interact with humans. Designing behaviors for social robots is a task that integrates knowledge form multiple scientific areas where knowledge about environment is paramount. This paper presents preliminary ideas on the modeling of social environment that are suitable to control purposes.

Social-intelligent robots should have four components: "act in ... complicated domains; communicate with humans using a language-like modality; reason about its actions at some level so that it has something to discuss, and, learn and adapt ... on the basis of human feedback", (Lopes and Connell, 2001).

The first component is already implemented on the MOnarCH robot in the sense that it is able to co-exist with people at a hospital ward and perform some tasks such as move around and occasionally say some sentences. The second and third components require perception for which this work provides a basic component. The goal is to build a motion awareness system for the robot providing spatial perception of the immediate surroundings and improving the level of interaction with people in the ward.

Figure 1 shows the MOnarCH robot used in this work. It was built to integrate social environments, namely the Pediatrics ward of a hospital. The robot interacts with children (the inpatients), and adults (the relatives and staff). As recognized, for instance, in (Holzinger et al., 2008) non-lab environments pose hard challenges in what concerns the acceptance of social robots due to the, usually high, number of factors involved. This means that perception information, namely from the neighbourhood of the robot is paramount. A Kinect sensor mounted on the head is used to provide range information from which models of the dynamics of the environment are created. Its actions will depend on these dynamics.



Figure 1: The MOnarCH robot.

The paper is divided as follows. Section 3 describes the motion awareness system and section 4 tests some aspects of the Kinect performance. An overview of the robot behavior implemented is explained in section 4.3. Conclusions and future work are discussed in section 5.

## 2 RELATED WORK

Sensors as Microsoft Kinect and Asus Xtion have been tested in robots to extract RGB-D images and

build 3D point-cloud representations of the environments, (Oliver et al., 2012). Both sensors proved to be suitable for mobile robotic navigation despite some limitations: narrowed field of view, small ranges, and measurements accuracy, (Eriksson and Ragnerius, 2012).

The Kinect has been shown to be useful in different areas such as robotics, performing arts, education, retail services and security, (Lun and Zhao, 2015), in 3D reconstructions of environments, (Zhang et al., 2015), and objects, (Varanasi and Devu, 2016), human movement recognition, (Cippitelli et al., 2016; Luo et al., 2014), and navigation with emphasis on obstacle avoidance, (Correa et al., 2012). A model based approach to detect humans using a 2D head contour model and a 3D head surface model is described in (Xia et al., 2011). An algorithm for pedestrian contours detection, by merging RGB and Depth images from a Kinect, was developed in (Chen et al., 2016). A method for tracking individuals is also proposed in (Yang et al., 2016), consisting of a subtraction method for background frames of depth images. The goal was to detect people who are on the verge of falling. An algorithm for human action recognition exploiting the skeleton provided by a Kinect is described in (Cippitelli et al., 2016).

Social robots are autonomous machines designed to interact with humans and show social behaviors (KPMG, 2016). When creating those behaviors, people will tend to project their thoughts and behaviors on robots (Duffy, 2003), and try to make them behave as humans. A relevant question is: how can researchers create models that give robots social skills to change people's perception towards the robots? A few tried to find solutions to this question. What people think about the robot's movement as it follows behind a person has been studied in (Gockley et al., 2007). Two approaches are implemented (direction-following and path-following). Results showed that the first is more similar to how humans behave. A robot that stands in line with humans is studied in (Nakauchi and Simmons, 2000), on the assumption that for a robot to be social it should recognize and react to people's social actions.

## 3 MOTION AWARENESS

Individuals adjust their movements in response to their neighbors' movements and positions, (Herbert-Read, 2016). This is called Herd Behavior and is characteristic of both animals and humans, e.g., flocking birds, dolphins, and nest building ants. Stock market bubbles, crowds and everyday decision-making ex-

emplifies how this concept is present in humans societies and has been studied by philosophers such as S. Kierkegaard and F. Nietzsche. Wilfred Trotter popularized the topic applied to humans, (Trotter, 1916). Human behavior, namely in social groups, is influenced by each others actions, (Musse and Thalmann, 1997).

These biological features inspired the modeling of social environments from specific events. Proper identification of the dynamics of a social environment is an enabler for strategies of adjustment of behaviors to the environment and hence maximize acceptance.

The Kinect sensor provides RGB-D images at 30 Hz frame rate and 640 x 480 pixels resolution, suitable for indoors robotics applications. The depth component is used to detect motion as it is more efficient in the subtraction process between two depth images (Greff et al., 2012) considered (see ahead).

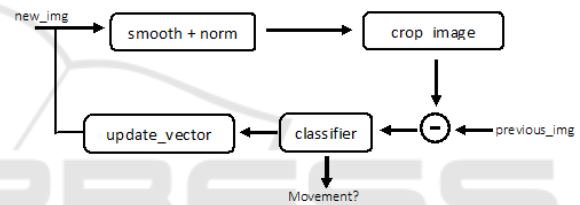Figure 2 shows the architecture of the motion awareness classifier developed.



Figure 2: Motion Awareness System.

The system starts by normalizing and smoothing each new image acquired by the Kinect. Normalization is done to fill in the pixels for which it was not possible for the Kinect to return any value, e.g., as when facing metallic areas. The smoothing process was done using the openCV function *medianBlur()*, with an linear aperture size of 5.

The images acquired from the Kinect are subject to noise caused by inadequate calibration, lighting conditions, and imaging geometry, and properties of the object surface, (Khoshelham and Elberink, 2012). This noise can be estimated under benign conditions, i.e., if the sensor is static and the environment conditions are constant, from the difference between the depth images obtained in close time intervals. Each depth image is mathematically represented as a matrix. The amount of noise can be estimated from the mean of all the depth values in the difference matrix. In lab conditions an average of 0.1 (normalized units) was found, which is the baseline value for non-lab trials. A noise free sensor would yield a zero image.

To reduce the computational cost only the bottom part of the image provided by the Kinect is used (*cut_image* function). Essentially, the goal is to determine if there are people passing by in the neigh-

borhood of the robot. As a person moves away from the sensor the downside part of the images still contains motion relevant information whereas the upper part may not contain any useful information (eventually related with the ceiling, in indoor environments). At the end of the process the new image has a resolution size of 640 x 380 pixels, smaller than the original resolution.

Figure 3) shows an example. In the top image no one is present. In the middle image someone is moving in front of the Kinect. The subtraction of the two images is shown at the bottom with the black regions representing unchanged areas (not affected by the movement) and the gray areas representing changes due to motion that occurred in the environment.


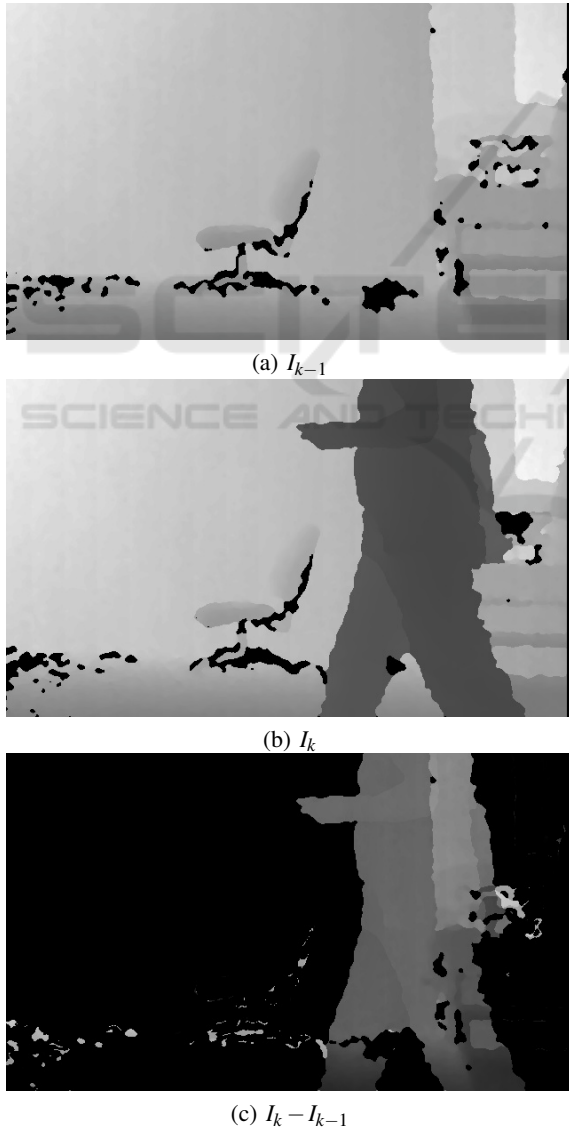(a) $I_{k-1}$


(b) $I_k$


(c) $I_k - I_{k-1}$

Figure 3: Example of depth images taken in the lab.

To classify the data obtained from the Kinect the estimates on the level of changes are compared with short term information (see Algorithm 1).

---

Algorithm 1: Classifier algorithm.

$N = 5$ {Short term memory window size}

$k_1 = 1$ {Time index, at 2.09 Hz average cycle rate}

$k_2 = 1$ {Short term memory update rate at 0.6Hz}

$T = 2$ {Decision threshold}

**Require:** Initialization procedure

**loop**
  Acquire $I_k$ {The depth image acquired at instant $k_1$}

  $d_{k_1} = \text{mean}(I_{k_1} - I_{k_1-1})$ {Mean differential depth}

  $D_{k_1} = [d_{k_1-1}, d_{k_1-2}, \ldots, d_{k_1-N}]$ {Short term memory of differential depths}

  $\overline{D}_{k_1} = \text{mean}(D_{k_1})$

  **if** Time to update short time memory **then**
    $\mathcal{D} = \left\{ |d_{k_1-j} - \overline{D}_{k_1}|, \; j = 0\ldots N \right\}$ {Short term deviations to the mean}

    $\mathcal{D}_s = \text{sort}(\mathcal{D})$ {Sorting in ascending order}

    $D_{k_1+1} = \mathcal{D}_s[1:N]$ {Keep the $N$ smallest deviations}

  **end if**
  **if** $|d_{k_1} - \overline{D}_{k_1}| > T$ **then**
    Signal movement detected
  **end if**

**end loop**

---

Algorithm 1 embeds a short term memory that allows the adaption of the system to changing environments. Even if no movement is detected the variable $\overline{D}_{k_1}$ keeps being updated. The speed of adjustment can be tuned through the short term window size, the update rate of $D_{k_1}$, and/or the decision threshold.

People walks at an average frequency of 1.9 Hz (normal pace) and average velocity of 1.35 m/s, (Ji et al., 2005). This means that a person can easily be caught in the Kinect field of view more than once as it moves in front of the sensor as the cycle rate of the classifier is higher. As the velocity decreases more times the person is detected and hence the detection events obtained does not discriminate among persons.

If necessary, this bias can be removed by imposing a minimum time between events of $1/2.09$ s.

Figure 4 shows $d_k$ (blue points) and $|d_k - \overline{D}_k|$ (black points) over time. The black points indicate instants where the classifier detected someone was nearby. This sample was obtained at the hospital with the sensor facing a lobby area (see figure 5). This location is a crossing point and usually there is people passing by. The detection of movement in a lobby area is a key feature for a social robot making it aware of the social space and thus enabling interesting interactions with people. Furthermore, people passing by may also be increasingly aware that the robot knows (in social terms) its surrounding space.
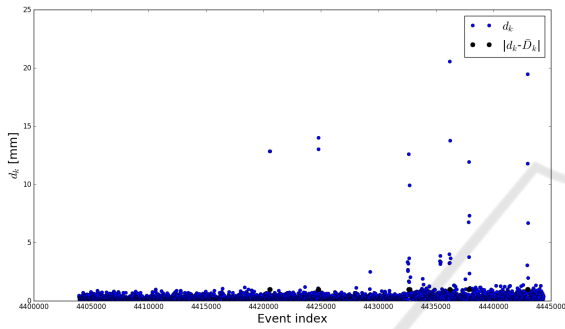

Figure 4: Classifier variables along time.


Figure 5: Hospital ward view from the sensor perspective.

The $d_k$ values around 0 represent moments of inactivity. Peaks correspond to movement detection (including false positives). The $\overline{D}_k$ tends to 0.2.

The short term memory is initialized to $D_0 = [0, 0, 0, 0, 0]$. As images are acquired, these values are adjusted (see algorithm 1).

As aforementioned, the Kinect produces some noise which can skew the $d_k$, hence the black points may refer to false positives. In this sample a 96% true positive was achieved, which is sufficient for HRI applications.

Figure 6 corresponds to the adjustment of the $\overline{D}_k$. This variable allows the system to infer if there is movement on the robot surroundings. At initialization $\overline{D}_k = 0.6$ because the values previously measured
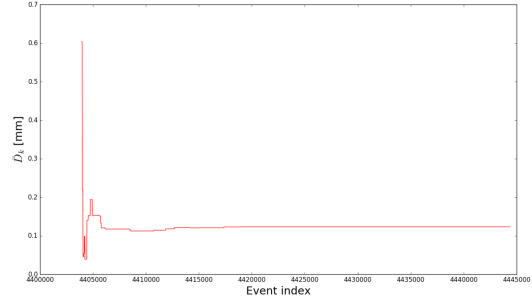

Figure 6: $\overline{D}_k$ temporal evolution.

and stored by the system were close to it. After the Kinect starts to adjust this value $\overline{D}_k$ starts to approximate 0.12, setting the value from which the system decides if reacts or no.

# 4 EXPERIMENTS

The experiments in this section illustrate the identification of statistical models for the social environment at the hospital from a perspective of a static bystander (in the case the robot with a Kinect sensor on board).

## 4.1 Setting up the System

The purpose of the initial tests is to verify that no exogenous conditions, e.g., lighting are likely to disturb the sensing (though they may still influence the data acquired), and establish a baseline for future work.

### 4.1.1 Test 1

Figure 7 shows images taken at the hospital where it can be seen that two people are staying in front of the robot. The upper row shows the previous and current images. In the bottom row, the first one shows the difference between previous image and the current one whereas the second shows the same image obtained immediately after the motion (and hence containing only the detection noise).

The second image on the bottom row shows the difference if people were not present, i.e., the image after the subtraction process. The $d_k$ value of each image on the bottom row, are 4.5 and 0.8, respectively. The $\overline{D}_k$ at this point was 0.45. Applying the threshold on this information one concludes that on the first situation people were present and on the second one were not.
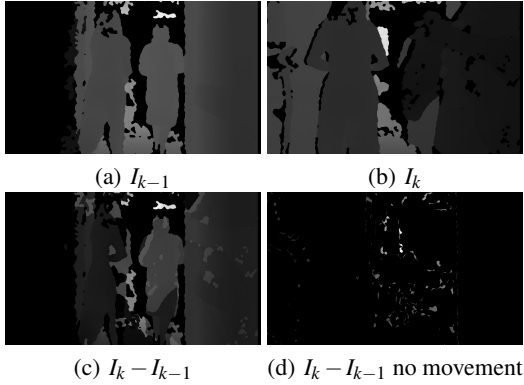
(a) $I_{k-1}$     (b) $I_k$

(c) $I_k - I_{k-1}$     (d) $I_k - I_{k-1}$ no movement

Figure 7: Experiment 1 (hospital).

### 4.1.2 Test 2

The data acquired by the Kinect is not always accurate, with the images showing some noise (see the comments in section 3). Unlike RGB images, depth images should be less influenced by lighting conditions. Nevertheless, if the Kinect is outshine by direct light, depth measurements are affected as these images are created using the infrared laser that uses infrared patterns for depth estimation. Thus infrared patterns provide wrong measurements if blinded by light. This results in undefined pixels on those areas. If the lighting changes, the undefined pixels might also change to specific values, or remain undefined.



(a) $I_{k-1}$     (b) $I_k$

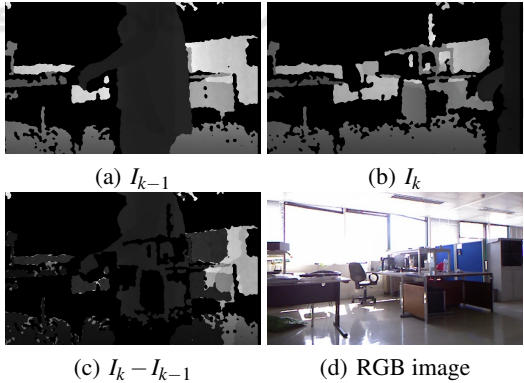(c) $I_k - I_{k-1}$     (d) RGB image

Figure 8: Experiment 2 (lab).

To ensure that the implementation is robust to lighting changes the robot was placed at a specific location inside the lab, facing a window. The RGB image of this example is shown on the second image in the bottom row in figure 8.

Figure 8 also shows the previous and current depth images at a time when someone was passing by, and the correspondent image difference. When the images were taken $\overline{D}_k = 0.103$ and $d_k = 8.29$. The classifier concluded that there was someone nearby.

Comparing these values with those of the first test it can be argued that there are no significant differences in the $\overline{D}_k$ and $d_k$ values. The locations where both experiments were performed are different thus leading to different values. Even if the places were the same, these values would not be exactly equal because they are influenced by the distances to the obstacles in front of the robot.

In conclusion, the illumination does not impacts the decision making of whether the environment is dynamic or not.

## 4.2 Identification of the Hospital Environment

These experiments took place at the Pediatrics ward of a hospital, with the robot placed in a small lobby connecting the main corridor and a playroom for the inpatient children. The robot stayed static during the full trials period, with the Kinect sensor active between $8:00$ and $22:00$. This period was divided into a set of 5 smaller periods, of unequal length, empirically defined (see below).

Figure 9 plots the actual time between events along a period of 6 days. The raw event index plot does not classifies the events according to their occurrence. Instead it shows them in the occurring sequence. Both plots clearly show differences in activity though in the raw event index plot the regions of higher and lower activity are better discriminated.

Figure 10 shows a power law fitting to the time indexed data, with the corresponding parameters shown in Table 1 for the day periods considered.

Table 1: Parameters $a, b, c$ estimated for a power law of the form $a(x+b)^{-n} + c$, for $n = 2$.

| Periods (hours) | a | b | c |
|---|---|---|---|
| 8-11 | 1.802e4 | 1.785 | 21.05 |
| 11-14 | 1.643e4 | 3.028 | 8.464 |
| 14-16 | 1.674e4 | 1.64 | 12.5 |
| 16-20 | 2.024e4 | 1.417 | 12.63 |
| 20-22 | 1.9664 | 1.368 | 28.77 |

The curves are clearly comparable as their shape parameters are close to each others among the different periods. This supports the claims that (i) a social environment may be given a power law representation (for the time between people passing at a specific place), and (ii) the vast majority of the events exhibits a small time between them, expressed by the long flat plateau that follows event index 1000.

Moreover, figure 10 can be interpreted as a 2-state environment dynamics. If events are identified with the time between detections, a first state indicates a

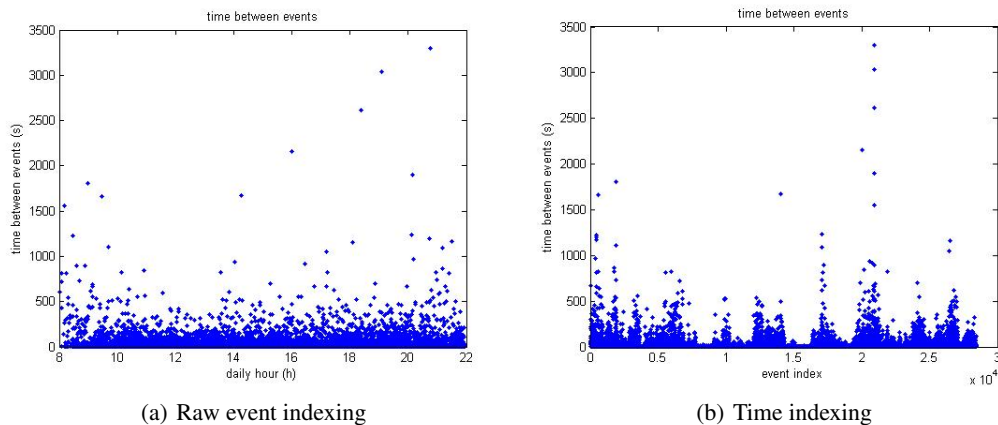(a) Raw event indexing            (b) Time indexing
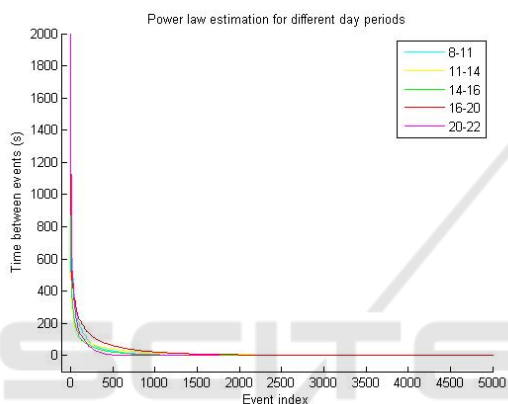
Figure 9: Time between events (s).



Figure 10: Power law fitting to time indexed data for the different daylight periods.

relatively small number (compared to the duration of the experiments) of long duration events. This corresponds to the left hand side of the plots, until around index 500. A second state indicates a big number (more frequent) of short duration events, roughly corresponding to the region after index 500.

Table 2 shows the parameters for Weibull, Gamma, and Lognormal distributions estimated over the aforementioned periods. These distributions were selected as (i) they cover the positive numerical space, and (ii) they are parameterized such that comparable shapes can be obtained by manipulation of their parameters.

The numerical values suggest some accordance between Weibull and Lognormal distributions. These can thus be used to represent the arousal conditions of the social environment and establish activity levels compatible with them.

## 4.3 Towards Robot Behaviors in a Social Environment

How a robot should behave in a social environment and what people expect from it is still a trendy topic in human-robot interaction. A common objective among humans is to obtain social acceptance. In fact this is at the core of human personality (see for instance (Maslow, 1970)). In what concerns social robotics, this is also a natural objective though it is entirely admissible that a robot be introduced in a social environment with disruptive goals (something that also happens in human environments).

The behaviors of the robot must make people believe that the robot is aware of the daily dynamics of the environment and that it can adapt itself and foster its own social integration, contributing to improve the global mood and offer alternatives to the people's routine.

The statistical models in the previous section suggest that some behaviors be designed using also a statistical approach. For example, having the robot exhibiting liveliness features is likely to contribute to the acceptance by the people present in the ward. However, these must match corresponding features of the environment, this meaning that some statistical models are not to be disturbed.

A single behavior is used to assess the performance of the classifier, namely implementing a liveliness feature. The robot simply turns the head at regularly spaced intervals, every 4 minutes on average. The goal is to assess the changes in the models previously identified.

For staff the presence of a robot should not interfere with the normal operation of the ward. This means that small or no disturbances were detected in the environment models identified. Inpatients and visitors are more likely to stay in front of the robot,

Table 2: First and second moments, $\mu, \sigma^2$, per distribution and period considered.

| Distribution | Periods (hours) | | | | |
|---|---|---|---|---|---|
| | 8-11 | 11-14 | 14-16 | 16-20 | 20-22 |
| | $\mu, \sigma^2$ | $\mu, \sigma^2$ | $\mu, \sigma^2$ | $\mu, \sigma^2$ | $\mu, \sigma^2$ |
| Weibull | 15.26, 1.35e3 | 6.81, 157.48 | 10.12, 436.02 | 9.01, 380.87 | 22.56, 3.995e3 |
| Gamma | 23.03, 1.68e3 | 8.92, 182.98 | 14.31, 547.98 | 13.70, 537.51 | 34.63, 4.341e3 |
| Lognormal | 11.79, 2.64e3 | 5.04, 134.46 | 7.58, 508.61 | 6.28, 349.54 | 18.88, 1.634e4 |

Table 3: First and second moments, $\mu, \sigma^2$, per distribution and period considered with liveliness behavior.

| Distribution | Periods (hours) | | | | |
|---|---|---|---|---|---|
| | 8-11 | 11-14 | 14-16 | 16-20 | 20-22 |
| | $\mu, \sigma^2$ | $\mu, \sigma^2$ | $\mu, \sigma^2$ | $\mu, \sigma^2$ | $\mu, \sigma^2$ |
| Weibull | 9.68, 344.22 | 4.52, 42.51 | 8.23, 195.08 | 5.46, 84.06 | 9.74, 384.73 |
| Gamma | 12.15, 350.14 | 5.21, 42.44 | 9.72, 190.05 | 7.03, 99.67 | 12.94, 424.92 |
| Lognormal | 7.58, 434.94 | 3.43, 22.65 | 6.63, 220.76 | 3.8, 40.2 | 7.09, 378.96 |

leading to periods of significant decreases in the time between events. Moreover, if the head of the robot moves, the perception of liveliness should increase and this should be reflected in the models estimated both (i) because of people responding to the head movement, and (ii) the movement of the head itself that causes the depth image subtraction to have enough information for the classifier to decide for movement.

Table 3 shows the parameters estimated for the same three distributions considered before. These results were obtained for a single day, without the first period.

The results suggest that even simple primitive behaviors such as having the robot turning the head on a regular basis have an impact on the models. Whatever the distribution chosen, there is a decreasing trend in the mean values, namely in the three afternoon periods, as expected. The variances in these periods show a decrease trend and, overall, the results are consistent with the expectations.

These quantitative results suggest that very simple interactions may indeed induce changes in the social environment. Moreover, they provide an objective measure of acceptance. In fact, the number of events occurring within a certain interval may be interpreted as an acceptance indicator. One can have situations where the higher the number of events the greater the acceptance; if persons tend to move in front of a robot more frequently when it moves the head this may indicate acceptance of the motion (people tend naturally to avoid repeating non-rewarding behaviors).

## 5 CONCLUSIONS AND FUTURE WORK

The paper discussed the implementation of a motion awareness system and its use in the identification of models for social environments based on time between events.

The system is based in depth image information with a short term memory classifier.

Section 4 tested the system against exogenous conditions, that could affect how images are provided. The conclusion was that the normal lighting variations do not influence the system. The results obtained also in this section show that a hospital ward environment can be modeled by a deterministic power law and by probability distributions, namely Weibull and Lognormal. These models are enablers to further studies on (i) other classifiers, and (ii) adjustment strategies for the behaviors of the robot.

Future work extends this discussion to (i) a moving sensor, as when the robot moves along the ward (and not only to the rotation of the head), and (ii) specific classes of events, e.g., events observed in specific areas of the environment or events observed from depth and color information.

Also, acceptance has been shown to depend on the educational background of people. By combining Kinect based data from a system, such as the one described in this paper, with a people identification strategy it will be possible to identify which groups of people influence acceptance (recall that in a telemedicine context (Ziefle et al., 2013) have shown that different social groups accept technologies differently).

# ACKNOWLEDGEMENTS

# REFERENCES

Chen, X., Henrickson, K., and Wang, Y. (2016). Kinect-based pedestrian detection for crowded scenes. *Computer-Aided Civil and Infrastructure Engineering*, 31(3):229–240.

Cippitelli, E., Gasparrini, S., Gambi, E., and Spinsante, S. (2016). A human activity recognition system using skeleton data from rgbd sensors. *Computational intelligence and neuroscience*, 2016:21.

Correa, D. S. O., Sciotti, D. F., Prado, M. G., Sales, D. O., Wolf, D. F., and Osorio, F. S. (2012). Mobile robots navigation in indoor environments using kinect sensor. In *Critical Embedded Systems (CBSEC), 2012 Second Brazilian Conference on*, pages 36–41. IEEE.

Duffy, B. (2003). Anthropomorphism and the social robot. *Robotics and autonomous systems*, 42(3):177–190.

Eriksson, T. and Ragnerius, E. (2012). Autonomous robot navigation using the utility function method and microsoft kinect.

Gockley, R., Forlizzi, J., and Simmons, R. (2007). Natural person-following behavior for social robots. In *Proceedings of the ACM/IEEE international conference on Human-robot interaction*, pages 17–24. ACM.

Greff, K., Brandão, A., Krauß, S., Stricker, D., and Clua, E. (2012). A comparison between background subtraction algorithms using a consumer depth camera. In *VISAPP (1)*, pages 431–436.

Herbert-Read, J. (2016). Understanding how animal groups achieve coordinated movement. *Journal of Experimental Biology*, 219(19):2971–2983.

Holzinger, A., Schaupp, K., and Eder-Halbedl, W. (2008). An Investigation on Acceptance of Ubiquitous Devices for the Elderly in an Geriatric Hospital Environment: using the Example of Person Tracking. In Miesenberger, K., Klaus, J., Zagler, W., and Karshmer, A., editors, *Computers Helping People with Special Needs*, Lecture Notes in Computer Science, LNCS 5105, pages 22–29. Heidelberg, Berlin: Springer.

Ji, T. et al. (2005). Frequency and velocity of people walking. *Structural Engineer*, 84(3):36–40.

Khoshelham, K. and Elberink, S. O. (2012). Accuracy and resolution of kinect depth data for indoor mapping applications. *Sensors*, 12(2):1437–1454.

KPMG (2016). *Social Robots*. KPMG.

Lopes, L. S. and Connell, J. H. (2001). Semisentient robots: Routes to integrated intelligence. *IEEE Intelligent Systems*, 16(5):10–14.

Lun, R. and Zhao, W. (2015). A survey of applications and human motion recognition with microsoft kinect. *International Journal of Pattern Recognition and Artificial Intelligence*, 29(05):1555008.

Luo, J., Wang, W., and Qi, H. (2014). Spatio-temporal feature extraction and representation for rgb-d human action recognition. *Pattern Recognition Letters*, 50:139–148.

Maslow, A. (1970). *Motivation and Personality*. Harper & Row, N.Y.

Musse, S. and Thalmann, D. (1997). A model of human crowd behavior: Group inter-relationship and collision detection analysis. In *Computer animation and simulation*, volume 97, pages 39–51. Springer.

Nakauchi, Y. and Simmons, R. (2000). A social robot that stands in line. In *Intelligent Robots and Systems, 2000.(IROS 2000). Proceedings. 2000 IEEE/RSJ International Conference on*, volume 1, pages 357–364. IEEE.

Oliver, A., Kang, S., Wünsche, B. C., and MacDonald, B. (2012). Using the kinect as a navigation sensor for mobile robotics. In *Proceedings of the 27th Conference on Image and Vision Computing New Zealand*, pages 509–514. ACM.

Trotter, W. (1916). *Instincts of the Herd in Peace and War*. Macmillan.

Varanasi, S. and Devu, V. K. (2016). 3d object reconstruction using xbox kinect v2. 0.

Xia, L., Chen, C.-C., and Aggarwal, J. K. (2011). Human detection using depth information by kinect. In *CVPR 2011 WORKSHOPS*, pages 15–22. IEEE.

Yang, L., Ren, Y., and Zhang, W. (2016). 3d depth image analysis for indoor fall detection of elderly people. *Digital Communications and Networks*, 2(1):24–34.

Zhang, J., Huang, Q., and Peng, X. (2015). 3d reconstruction of indoor environment using the kinect sensor. In *Instrumentation and Measurement, Computer, Communication and Control (IMCCC), 2015 Fifth International Conference on*, pages 538–541. IEEE.

Ziefle, M., Klack, L., Wilkowska, W., and Holzinger, A. (2013). Acceptance of Telemedical Treatments A Medical Professional Point of View. In Yamamoto, S., editor, *Human Interface and the Management of Information. Information and Interaction for Health, Safety, Mobility and Complex Environments*, Lecture Notes in Computer Science LNCS 8017, pages 325–334. Berlin Heidelberg: Springer.